

Project By: Hritul Pardhi 18070123061 ENTC - 'A'

Description: In this program of sentimental analysis of hotel reviews, I have taken the data of reviews of hotel, in which reviews had a column of happy or not happy, so based on the words used in the review and whether the person was happy or not happy after his visit at that hotel, gives me the data as an input in this project to feed to sklearn for training.

After the modelling and creating this project: One can write a review and check it automatically whether the person is happy or not happy after his/her visit.

```
1 from google.colab import drive
2 drive.mount('/gdrive')
```

Drive already mounted at /gdrive; to attempt to forcibly remount, call drive.mount("/gdr



```
1 import warnings
2 warnings.filterwarnings('ignore')
```

```
1 # Mounting Drive
2 import os
3 os.chdir('/gdrive/My Drive/')
```

Data Facts and Importing Libraries

```
1 import pandas as pd
```

```
1 df = pd.read_csv('train.csv')
```

```
1 df.shape

(38932, 5)
```

```
1 df.head()
```

```
User_ID      Description  Browser_Used  Device_Used  Is_Response

1 df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 38932 entries, 0 to 38931
Data columns (total 5 columns):
#   Column          Non-Null Count  Dtype
---  -
0   User_ID         38932 non-null  object
1   Description      38932 non-null  object
2   Browser_Used    38932 non-null  object
3   Device_Used     38932 non-null  object
4   Is_Response     38932 non-null  object
dtypes: object(5)
memory usage: 1.5+ MB
```

```
1 df.describe().transpose()
```

	count	unique	top	freq
User_ID	38932	38932	id23295	1
Description	38932	38932	Stayed for recent convention. Perfect location...	1
Browser_Used	38932	11	Firefox	7367
Device_Used	38932	3	Desktop	15026
Is_Response	38932	2	happy	26521

Data Cleaning / Data Analysis

```
1 ## Checking missing vales in data set and printing percentage of missing values
2 ## for each column
3
4 count = df.isnull().sum().sort_values(ascending=False)
5 percentage = ((df.isnull().sum()/len(df)*100)).sort_values(ascending=False)
6 missing_data = pd.concat([count,percentage],axis=1, keys=['Count','Percentage'])
7
8 print('Count and percentage of missing values for the columes:')
9 missing_data
```

Count and percentage of missing values for the columns:

```

Count Percentage
-----
1  ## No missing values found
   -----
1  ## checking for the distribution of default
2  import matplotlib.pyplot as plt
3  %matplotlib inline
4  print('Percentage for defalut\n')
5  print(round(df.Is_Response.value_counts(normalize=True)*100,2))
6  round(df.Is_Response.value_counts(normalize=True)*100,2).plot(kind='bar')
7  plt.title('Percentage Distributions by review type')
8  plt.show()

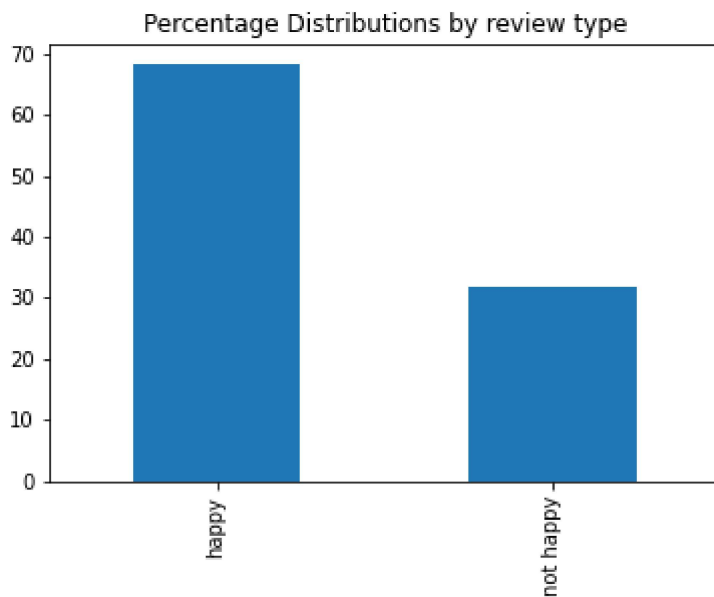
```

Percentage for defalut

```

happy      68.12
not happy   31.88
Name: Is_Response, dtype: float64

```



```

1  # Removing Extra columns
2  df.drop(columns=['User_ID', 'Browser_Used', 'Device_Used'], inplace=True)

1  df.head()

```

```

                                Description  Is_Response
0    The room was kind of clean but had a VERY stro...    not happy

1  # Applying cleaning for text
2  import re
3  import string
4  # Function converts to lower case, removes square brackets, removes numbers and
5  # punctuations
6  def text_clean_1(text):
7      text=text.lower() #lowercase
8      text=re.sub('\[.*?\]', '', text) #brackets
9      text=re.sub('[%s]' % re.escape(string.punctuation), '', text) #punctuations
10     text=re.sub('\w*\d\w*', '', text) #digits
11     text = re.sub('[\'\""...]', '', text) #quotation mark
12     text = re.sub('\n', '', text) #next line/blank lines
13     return text
14
15  cleaned1 = lambda x: text_clean_1(x)

1  df['cleaned_description'] = pd.DataFrame(df.Description.apply(cleaned1))
2  df.head()

```



	Description	Is_Response	cleaned_description
0	The room was kind of clean but had a VERY stro...	not happy	the room was kind of clean but had a very stro...
1	I stayed at the Crown Plaza April -- - April -...	not happy	i stayed at the crown plaza april april th...
2	I booked this hotel through Hotwire at the low...	not happy	i booked this hotel through hotwire at the low...
3	Stayed here with husband and sons on the	not happy	stayed here with husband and sons on

MODEL TRAINING

```

1  from sklearn.model_selection import train_test_split
2
3  Independent_var = df.cleaned_description
4  Dependent_var = df.Is_Response
5
6  IV_train, IV_test, DV_train, DV_test = train_test_split(Independent_var, Dependent_var,
7
8  print('IV_train :', len(IV_train))
9  print('IV_test  :', len(IV_test))
10 print('DV_train :', len(DV_train))
11 print('DV_test  :', len(DV_test))

IV_train : 35038
IV_test  : 3894

```

```
DV_train : 35038
DV_test  : 3894
```

```
1 from sklearn.feature_extraction.text import TfidfVectorizer #importing essential librari
2 from sklearn.linear_model import LogisticRegression
3
4 tvec = TfidfVectorizer()
5 clf2 = LogisticRegression(solver = "lbfgs")
6
7
8 from sklearn.pipeline import Pipeline #for step by step execution
```

```
1 model = Pipeline([('vectorizer',tvec),('classifier',clf2)])
2
3 model.fit(IV_train, DV_train)
4
5
6 from sklearn.metrics import confusion_matrix
7
8 predictions = model.predict(IV_test)
9
10 confusion_matrix(predictions, DV_test)
```

```
array([[2417,  304],
       [ 154, 1019]])
```

```
1 from sklearn.metrics import accuracy_score, precision_score, recall_score
2
3 print("Accuracy : ", accuracy_score(predictions, DV_test))
4 print("Precision : ", precision_score(predictions, DV_test, average = 'weighted'))
5 print("Recall : ", recall_score(predictions, DV_test, average = 'weighted'))
```

```
Accuracy :  0.8823831535695943
Precision :  0.8889271415963718
Recall :  0.8823831535695943
```

```
1 #Final output:
2 example = ["It was an wonderful experience, I would visit again"]
3 result = model.predict(example)
4
5 print(result)

['happy']
```

